

Fraud-NonFraud Detection

Download Data

Missing Value Imputation Technique Determination

Null percentage < 5%:  
Drop null value

5% < Null percentage < 30%:  
impute missing values using corresponding skills

Null percentage > 30%:  
Drop the whole column

Missing Value Imputation

PWD\_UPDT\_TS

Null percentage = 22.3%

It's a date value, and very hard to fill the null value, therefore, add a column called "missing\_PWD\_UPDT\_TS" indicate if PWD\_UPDT\_TS is missing. If missing, equals 1, else equals 0

CARR\_NAME, RGN\_NAME, STATE\_PRVNC\_TXT

Null percentage = 19.35%

It is categorical value and all three missing at the same time, therefore, add the null value to "unknown" category

DVC\_TYPE\_TXT & AUTHC\_SCNDRY\_STAT\_TXT

DVC\_TYPE\_TXT Null percent = 12.6%

I chose to use KNN classifier to impute the null value for DVC\_TYPE\_TXT and AUTHC\_SCNDRY\_STAT\_TXT because I found there are some association between ALERT\_TRGR\_CD, DVC\_TYPE\_TXT, AUTHC\_PRIM\_TYPE\_CD, and AUTHC\_SCNDRY\_STAT\_TXT

Label encoding  
KNN imputation  
Replace original null value with predicted value

CUST\_STATE

Null percentage = 0.26%

Directly drop the rows with null values

PH\_NUM\_UPDT\_TS

Null percentage = 50.6%

Drop the column

Handling Date&Time

Split the data according to "/" or "-", and add new year, month, day, hour, minute, second related column

Exploring Data Analysis

Numerical Data

Categorical Data

Date&Time data

Data Processing Conclusion

ACTN\_CD, ACTN\_INTNL\_TXT, TRAN\_TYPE\_CD only contain single value, will drop all of them

CUST\_ZIP and CUST\_STATE are perfectly correlated, will drop CUST\_ZIP

TRAN\_DT contains year, month, day of TRAN\_TS, but not hour, minute, second. Will drop TRAN\_DT

Drop ACTVY\_DT and all related columns

Drop PWD\_UPDT\_TS and all related columns

Drop CUST\_STATS

Drop PH\_NUM\_UPDT\_TS

Data Processing

Handling missing values as in Imputaion step

Split year, month, day, hour, minute, second of each datetime feature into different columns

Drop useless features as suggested in data processing conclusion

Change year, month, day, hour, minute, second features to numerical value

Change target value in training dataset to binary format for modeling

Encoding

Binary Encoding

CARR\_NAME and STATE\_PRVNC\_TXT have too many distinct categories. One-hot-encoding will generate extremely sparse data. Therefore, choose binary encoding

One-hot-encoding

Rest of categorical data will perform one-hot-encoding

Model Training & Selection

Training

Ridge&Lasso Regression

Logistic Regression

Dummy Classifier

KNN

SVM

Decision Tree

Random Forest

Gradient Boosting

AdaBoosting

XGBoosting

Selecting

Select XGBoost since it has best score

Parameter Tunning

learning\_rate

max\_depth

min\_child\_weight

subsample

colsample\_bytree

n\_estimators

Result Prediction